

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26

Global threat of arsenic in groundwater

Joel Podgorski^{1,2} and Michael Berg^{1,3*}*

1. Eawag, Swiss Federal Institute of Aquatic Science and Technology, Department Water Resources and Drinking Water, 8600 Dübendorf, Switzerland
2. University of Manchester, Department of Earth and Environmental Sciences, Manchester, United Kingdom
3. UNESCO Chair on Groundwater Arsenic within the 2030 Agenda for Sustainable Development and School of Civil Engineering and Surveying, University of Southern Queensland, 4350 QLD, Australia

**Corresponding authors. Email: joel.podgorski@eawag.ch, michael.berg@eawag.ch.*

Revision for “*Science*”

Abstract:

Naturally occurring arsenic in groundwater affects millions of people worldwide. We created a global prediction map of groundwater arsenic exceeding 10 µg/L using a random forest machine learning model based on eleven geospatial environmental parameters and over 50,000 aggregated data points of measured groundwater arsenic concentration. Our global prediction map includes known arsenic-affected areas and previously undocumented areas of concern. Combining the global arsenic prediction model with household groundwater-usage statistics, we estimate that 94-220 million people are potentially exposed to high arsenic concentrations in groundwater, the vast majority being in Asia (94%). Since groundwater is increasingly utilized to support growing populations and buffer against water scarcity due to changing climate, this work is important to raise awareness, identify areas for safe wells, and help prioritize testing.

27 The natural, or geogenic, occurrence of arsenic in groundwater is a global problem with wide-
28 ranging health effects for humans and wildlife. Being toxic and not serving any beneficial
29 metabolic function, inorganic arsenic (the species present in groundwater) can lead to
30 disorders of the skin and vascular and nervous systems as well as cancer(1, 2). The major
31 source of inorganic arsenic in the diet is through arsenic-contaminated water, although
32 ingestion through food, particularly rice, represents another important route of exposure(3).
33 As a consequence, the World Health Organization (WHO) has set a guideline concentration of
34 10 µg/L in drinking water(4).

35 At least trace amounts of arsenic occur in virtually all rocks and sediments around the
36 world(5). However, in most of the large-scale cases of geogenic arsenic contamination in
37 groundwater, arsenic accumulates in aquifers composed of recently deposited alluvial
38 sediments. Under anoxic conditions, arsenic is released from the microbial and/or chemical
39 reductive dissolution of arsenic-bearing iron(III) minerals in the aquifer sediments(6-9). Also
40 under oxidizing high-pH conditions, arsenic can desorb from iron and aluminum
41 hydroxides(10). Furthermore, aquifers in flat-lying sedimentary sequences generally have a
42 small hydraulic gradient, causing groundwater to flow slowly. This longer groundwater
43 residence time allows dissolved arsenic to accumulate and its concentration to increase. Other
44 processes responsible for arsenic release into groundwater include oxidation of arsenic-
45 bearing sulfide minerals as well as release from arsenic-enriched geothermal deposits.

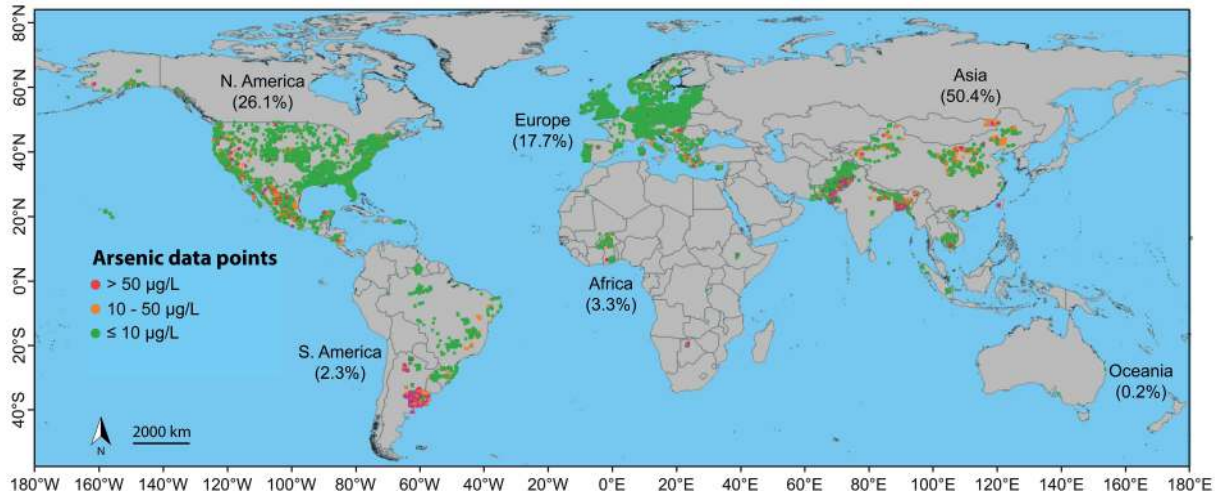
46 The fact that arsenic is generally not included in the standard suite of tested water quality
47 parameters(11) and is not detected by the human senses mean that arsenic is regularly being
48 discovered in new areas. Since one of the greatest occurrences of geogenic groundwater
49 arsenic was discovered in 1993 in the Bengal Delta (5, 12, 13), high arsenic concentrations
50 have been detected all around the world, with hotspots including Argentina(14-17),
51 Cambodia(18, 19), China(20-22), India(23-25), Mexico(26, 27), Pakistan(28, 29), the USA(30,
52 31) and Vietnam(32, 33).

53 In order to help identify areas likely to contain high concentrations of arsenic in groundwater,
54 several researchers have used statistical learning methods to create arsenic prediction maps
55 based on available datasets of measured arsenic concentrations and relevant geospatial
56 parameters. Previous studies have focused on Burkina Faso(34), China(21, 35), South Asia(29,
57 36), Southeast Asia(37), the USA(31, 38, 39) and the Red River Delta in Vietnam(33) as well as
58 sedimentary basins around the world (40). The predictor variables used in these studies
59 generally include various climate and soil parameters, geology and topography (Table S3).

60 Taking advantage of the increasing availability of high-resolution datasets of relevant
61 environmental parameters, we use statistical learning to model what to our knowledge is the
62 most spatially extensive compilation of arsenic measurements in groundwater assembled,
63 which makes a global model possible. In order to focus on health risks, we consider the
64 probability of arsenic in groundwater exceeding the WHO guideline. For this we have chosen
65 the random forest method, which our preliminary tests showed to be highly effective in
66 addressing this classification problem. We use the resulting model to produce the most

67 accurate and detailed global prediction map to date of geogenic groundwater arsenic, which
68 can be used to help identify previously unknown areas of arsenic contamination as well as
69 more clearly delineate the scope of this global problem and considerably increase awareness.

70



71

72 **Figure 1.** Arsenic concentrations excluding those known to originate from a depth greater
73 than 100 m. Values are from the sources listed in Table S1. The geographical distribution of
74 data is indicated by continent.

75

76

77 RESULTS

78 Random forest modeling

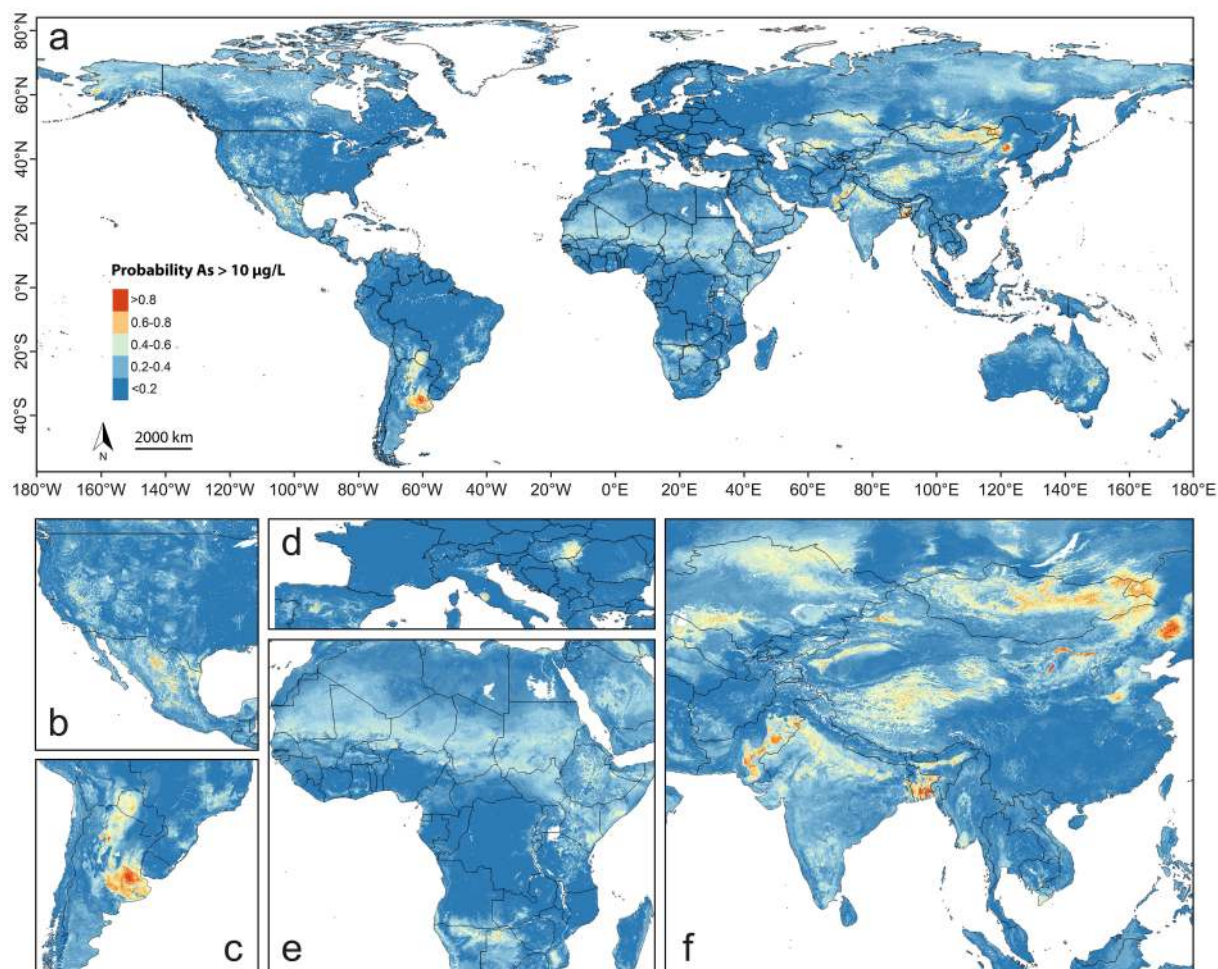
79 We aggregated data from nearly 80 studies of arsenic in groundwater (see Table S1 for
80 references and statistics) into a single dataset ($n > 200,000$). Averaging into 1 km^2 pixels
81 resulted in more than 55,000 arsenic data points for use in modeling based on groundwater
82 samples not known to originate from greater than 100 m depth (Figure 1).

83 In order to create the simplest and most accurate model, an initial set of 52 potentially
84 relevant environmental predictor variables was iteratively reduced in consideration of their
85 relative importance and impact on the accuracy of a succession of random forest models. The
86 final selection of 11 predictor variables (Table S2) includes several soil parameters (topsoil
87 clay, subsoil sand, pH and fluvisols), all of the climate variables (precipitation, actual and
88 potential evapotranspiration and combinations thereof as well as temperature) and the
89 topographic wetness index. In contrast, none of the geology variables proved to be
90 statistically important. This is not to imply that geology does not play a role in geogenic
91 arsenic accumulation, but rather that the particular geology variables tested were not as
92 relevant as the other variables. This may be due to the coarse nature of the geological maps,
93 which are standardized for the entire world. Although the number of predictor variables was

94 reduced by nearly 80%, both the Area Under the Curve (AUC, 0.89) and Cohen's kappa
95 statistic (0.55) remained unchanged.

96 The final random forest model was created based on the compiled global dataset of high and
97 low arsenic concentrations along with the 11 predictor variables. The standard number of
98 variables to be made available at each branch of each tree is between three and four (see
99 Methods). Since our tests showed the value of three performing better than four and higher
100 values (though error/performance rates varied only within ~1%), we set this parameter to
101 three. The global map produced from this model is displayed in Figure 2a along with more
102 detailed views of the more populated affected continental regions (Figure 2b-f). It indicates
103 the probability of the concentration of arsenic in groundwater in a given 1 km² cell exceeding
104 10 µg/L. The uncertainty of the model is inherent in the probabilities themselves, since they
105 are simply the average of the votes or predictions of high or low values of each of the 10,001
106 trees grown. That is, each tree casts a vote of 0 or 1 ('no' or 'yes' to As >10 µg/L) for each cell
107 based on the values of the predictor variables in that cell. Figures S2-S8 also provide more
108 detailed views of the prediction map for each of the inhabited continents.

109 The importance of each of the 11 predictor variables in terms of mean decrease in accuracy
110 and mean decrease in the Gini index is listed in Figure S1. Relative to the initial set of 52
111 variables, the values of these two statistics for most of the 11 final predictor variables appear
112 to fall within a fairly narrow range, indicating comparable importance. Exceptions include
113 fluvisols and soil pH, which have somewhat greater importance, and temperature, which
114 according to both statistics is the least important of the 11 variables. Soil pH was also found to
115 be an important predictor variable in arid, oxidizing environments in Pakistan(29). Although
116 widespread arsenic dissolution occurs in Holocene fluvial sediments(5-7, 9, 37), this geological
117 epoch has not been consistently mapped around the world. However, the global dataset of
118 fluvisols provides a very suitable alternative(29), which may even be more appropriate as
119 fluvisols by definition encompass recent fluvial sediments and not, for example, aeolian
120 Holocene sediments that are generally not relevant for arsenic release. The generally high
121 model importance of climate variables, as evidenced by them all being selected for the final
122 model, highlights the strong control that climate has on arsenic release in aquifers. In
123 particular, precipitation and evapotranspiration have a direct role in creating conditions
124 conducive for arsenic release under reducing conditions (e.g. waterlogged soils) as well as
125 high aridity associated with oxidizing high-pH conditions.



126

127 **Figure 2.** Modeled probability of arsenic concentration in groundwater exceeding 10 µg/L for
 128 the entire globe (a) along with zoomed-in sections of the main more densely populated
 129 affected area (b)-(f). The model is based on the arsenic data points in Figure 1 and the
 130 predictor variables in Table S2. Figures S2 to S8 provide more detailed views of the prediction
 131 map.

132

133 The performance of the random forest model on the test dataset (20% of the data, which was
 134 randomly selected while maintaining the relative distribution of high and low values) is
 135 summarized in the confusion matrix in Table 1. Despite a prevalence of high values (>10 µg/L)
 136 of only 22% in the dataset, the model performs well in predicting both high values (sensitivity:
 137 0.79) and low values (specificity: 0.85) at a probability cutoff of 0.50. The average of these two
 138 figures, known as balanced accuracy, is correspondingly high at 0.82. Likewise, the model's
 139 AUC, which considers the full range of possible cutoffs, has a very high value of 0.89 with the
 140 test dataset (Table 1). For comparison, the AUC of a random forest using all 52 original
 141 predictor variables is also 0.89.

142 The model was also tested on a dataset of over 49,000 arsenic data points originating from
 143 known depths greater than 100 m (average 562 m, standard deviation 623 m). Although the
 144 model was not trained on any measurements from these depths and the fact that only surface

145 parameters were used as predictor variables, the model nevertheless performed quite well in
 146 predicting the arsenic concentrations of these deep groundwater sources, as evidenced by an
 147 AUC of 0.77.

148

149

150 **Table 1.** Confusion matrix and other statistics summarizing the results of applying the random
 151 forest model to the test dataset at a probability cutoff of 0.50.

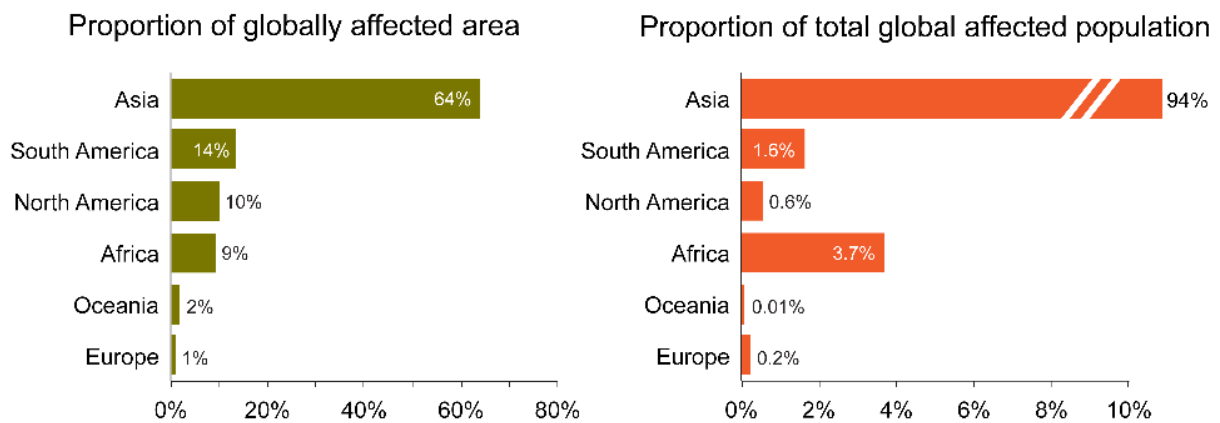
| | Measured | |
|---------------------------------|--------------|-------------------------|
| | As ≤ 10 µg/L | As > 10 µg/L |
| Predicted As ≤ 10 µg/L | 7710 | 555 |
| Predicted As > 10 µg/L | 1394 | 2037 |
| Sensitivity: 0.79 | | Prevalence: 0.22 |
| Specificity: 0.85 | | Balanced Accuracy: 0.82 |
| Positive Predictive Value: 0.59 | | Cohen’s kappa: 0.55 |
| Negative Predictive Value: 0.93 | | AUC: 0.89 |

152

153

154

155



156

157 **Figure 3.** Proportions of land area and population potentially affected by arsenic concentrations
 158 in groundwater exceeding 10 µg/L by continent.

159

160

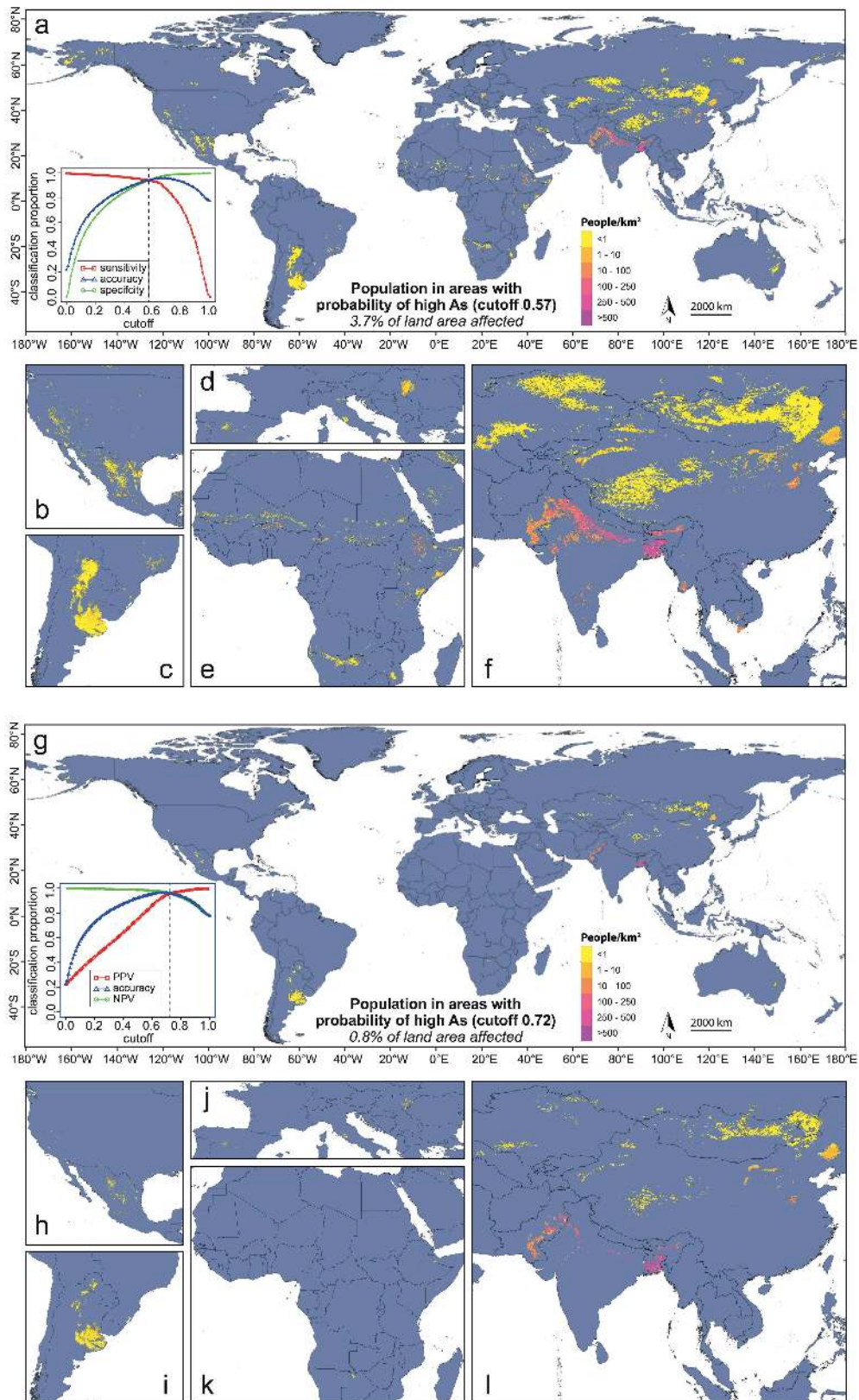
161 **Regions and populations at risk**

162 Areas predicted to have high arsenic concentrations in groundwater exist on all continents,
163 with most being located in Central, South and Southeast Asia, parts of Africa and North and
164 South America (Fig. 2 and Figs. S2-S8). Known areas of groundwater arsenic contamination are
165 generally well captured by the global arsenic prediction map, e.g. parts of the western US,
166 central Mexico, Argentina, the Pannonian Basin, Inner Mongolia, the Indus Valley, the Ganges-
167 Brahmaputra Delta and the Mekong River and Red River Deltas. Areas of elevated arsenic
168 hazard where little concentration data exist include parts of Central Asia, particularly
169 Kazakhstan, Mongolia and Uzbekistan, the Sahel region and broad areas of the Arctic and sub-
170 Arctic. Of these, the Central Asian hazard areas are better constrained as evidenced by higher
171 probabilities.

172 Probability threshold values of 0.57 from the sensitivity-specificity comparison and 0.72 from
173 the PPV-NPV comparison were found using the full dataset (combined training and test
174 datasets) of arsenic concentrations. The proportions of high modeled arsenic hazard by
175 continent associated with each of these probabilities are shown in Figure 3. Global maps of
176 the potentially affected population in the risk areas as determined by these two thresholds
177 are shown in Figure 4. As described in the Methods, these maps were then used to estimate
178 the population potentially affected by drinking groundwater with arsenic concentrations
179 exceeding 10 µg/L.

180 The resulting global arsenic risk assessment indicates that approximately 94-220 million
181 people around the world (of which 85-90% are in South Asia) are potentially exposed to high
182 concentrations of arsenic in groundwater from their domestic water supply (Tables S4 and
183 S5). This range is consistent with the previous most comprehensive literature compilations,
184 that is 140 million people(41) and 225 million people(42). Household groundwater-use
185 statistics were not available for ~6-8 percent of the affected countries (depending on the
186 cutoff), for which the less detailed statistics derived from the FAO AQUASTAT database were
187 used instead (see Methods for details). To determine the amount of error that using these
188 more general groundwater-use statistics might introduce to the overall population figures, the
189 global potentially affected populations were recalculated with these countries' (those lacking
190 household groundwater-use statistics) groundwater-use rates set to the extreme values of 0%
191 and 100%. Since this applied to relatively few countries and As-affected areas, doing so
192 affected the overall global population figures by an inconsequential amount ($\pm 0.1\%$),
193 indicating that using the AQUASTAT groundwater-use rates, where necessary, is an acceptable
194 approximation.

195



196
 197 **Figure 4:** Population in risk areas potentially containing aquifers with arsenic concentrations
 198 $>10 \mu\text{g/L}$ using probability cutoffs of (a) 0.57, at which sensitivity and specificity are equal
 199 (inset) as applied to the full (training and test) dataset, and (g) 0.72, at which PPV and NPV are
 200 equal (inset) using the full dataset. The detailed areas of Fig. 2 are also repeated here for both
 201 models (b)-(f) and (h)-(l).

202 This estimate of risk takes into account only the proportion of households utilizing
203 unprocessed groundwater and assumes uniform rates throughout the urban and non-urban
204 areas of each country. The uncertainties of these rates are unknown. The population in each
205 cell was reduced by the uncertainty of the cell's prediction, which is justified based on the
206 heterogeneity inherent in the accumulation of arsenic in an aquifer, which is generally at a
207 much finer scale than that of the 1-km² resolution of the arsenic hazard map. Since the arsenic
208 prediction for a cell represents the average outcome for that cell, we can take the modeled
209 probability as a first-order approximation of the proportion of an aquifer in that cell
210 containing high arsenic concentrations. Only cells exceeding the probability threshold (i.e.
211 0.57 or 0.72) were considered. The global estimate of 94-220 million people potentially
212 affected by consuming arsenic-contaminated groundwater is broken down by continent and
213 country in Tables S4 and S5, respectively, and represents the most accurate and consistent
214 global estimate available.

215 **DISCUSSION**

216 The accuracy of the global groundwater arsenic prediction model presented here, as
217 indicated, for example, with an AUC of 0.89 calculated with the test dataset, exceeds that
218 found in previous arsenic prediction studies (Table S3). The dominance of climate and soil
219 parameters in the final model is indicative of their direct influence or at least strong
220 association with the processes of arsenic accumulation in groundwater.

221 With respect to previous arsenic prediction maps of global sedimentary basins(40, 43), the
222 new model represents a significant advancement on a few different levels. First of all, the new
223 model presented here provides predictions for all areas of the inhabited continents, whereas
224 the previous first-generation statistical model covered only about half of the land areas. In
225 addition, a ten-fold increase in measurement points has allowed arsenic concentrations to be
226 incorporated from many more areas of the globe. The greatly expanded availability and
227 quality of global predictor datasets over the past ten years has enabled new variables to be
228 considered, such as soil type (e.g. fluvisols), as well as provided a 10 to 60-fold greater spatial
229 resolution (i.e. 30 arc-seconds versus 5-30 arc-minutes). However, the presence of high
230 arsenic in groundwater at a given location is of course predicated on the existence of an
231 aquifer in the first place, which may not be so in the case of unfractured solid rock, steep
232 terrain or very dry conditions. Models are only as good as the data on which they are based.
233 As accurate as the new arsenic model is, it could be further improved as more arsenic data
234 and more detailed predictor datasets come into existence.

235 Particularly in sedimentary aquifers, arsenic concentration is often highly dependent on
236 depth, that is, to specific sedimentary sequences that differ in the concentration of arsenic in
237 sediments as well as the geochemical conditions conducive to arsenic release. In order to
238 better characterize this relationship in a given sedimentary basin, detailed depth information
239 of groundwater samples would need to be incorporated in a separate basin-level study.
240 Unfortunately, it is not feasible in a global-scale study to account for all of the diversity of the
241 sedimentary basins of the world, especially since depth information of groundwater samples

242 is often not available. As such, we have relied on a statistical analysis of model performance
243 against depth ranges of samples (where present) in order to determine model sensitivity to
244 depth.

245 Our approach in the risk assessment of potentially affected population is relatively discerning
246 and/or conservative. As such, the resulting population estimates may in some cases be lower
247 than those found in earlier studies. One reason for this is that we employed country-specific
248 statistics of rural and urban domestic groundwater usage, which allowed us to subtract the
249 proportion of population utilizing surface water, tap water or other sources. This was not the
250 case, for example, in a previous study of China that estimated 19.6 million people being
251 affected in the country(21), whereas our estimate is considerably lower at 4.3-12.1 million.
252 Furthermore, we consider only areas in which the probability of high arsenic exceeds the
253 statistically determined cutoffs, i.e. 0.57 and 0.72. Taking the USA as an example, applying this
254 criterion left only 0.2-2% of the area of the country over which to sum the potentially affected
255 population (≤ 0.21 million, this study). In a previous arsenic risk assessment of the USA(31), the
256 entire country was used to estimate affected population (2.1 million), that is, not only the
257 high risk areas.

258 The actual proportion of groundwater usage varies spatially throughout a country, and so
259 more detailed usage statistics beyond only urban versus rural would improve the accuracy of
260 a risk assessment. In addition, more groundwater samples (ideally including depth
261 information) from areas that currently have poor coverage would benefit future modeling
262 efforts by allowing the model to be better adapted to those areas.

263 The presented arsenic probability maps should be used as a guide to further groundwater
264 arsenic testing, for example in Central Asia, the Sahel and other regions of Africa. Only actual
265 groundwater quality testing can definitively determine the suitability of groundwater with
266 respect to arsenic, particularly due to small-scale (<1 km) aquifer heterogeneities that cannot
267 be modeled with existing global datasets(9, 44). The hazard maps highlight areas at risk and
268 provide a basis for targeted surveys, which continue to be important. The already large
269 number of people potentially affected can be expected to increase as groundwater use
270 expands with a growing population and increasing irrigation, especially in the light of water
271 scarcity associated with warmer and drier conditions related to climate change. The maps can
272 also help aid mitigation measures, such as awareness raising, coordination of government and
273 financial support, health intervention programs, securing alternative drinking water resources
274 and arsenic removal options tailored to the local groundwater conditions as well as social
275 setting.

276

277 **ACKNOWLEDGEMENTS**

278 We thank our colleagues Anja Bretzler and Chris Zurbrügg (Eawag) and Andreas Steiner and
279 Stephanie Piers de Raveschoot (SDC), and D.A. Polya and R. Wu (University of Manchester) for
280 their support as well as the many providers of data, which were an essential component of

281 this work. **Funding:** We thank the Swiss Agency for Development and Cooperation (projects
282 no. 7F-09010.01.01 and 7F-09963.01.01) for their long-term support and co-funding of this
283 study, as well as a University of Manchester EPSRC IAA Impact Support Fund Award. **Author**
284 **contributions:** J.P.: Methodology, Modeling, Writing- Original draft preparation, M.B.:
285 Supervision, Writing- Reviewing and Editing. **Competing interests:** Authors declare no
286 competing interests. **Data and materials availability:** The modeling data, code and raster
287 output maps are available at ERIC/open(45). Arsenic concentration data points and hazard
288 and risk maps are also available for viewing on the GIS-based Groundwater Assessment
289 Platform (GAP), www.gapmaps.org.

290

291 **Supplementary Materials**

292 Methods

293 Supplementary Tables 1-6

294 Supplementary Figures 1-11

295 References 46-125

296

297

298 **References and notes**

- 299 1. A. H. Smith, E. O. Lingas, M. Rahman, Contamination of drinking-water by arsenic in
300 Bangladesh: a public health emergency. *Bulletin of the World Health Organization* **78**, 1093-
301 1103 (2000).
- 302 2. M. F. Hughes, Arsenic toxicity and potential mechanisms of action. *Toxicology letters* **133**, 1-16
303 (2002).
- 304 3. D. Mondal *et al.*, Comparison of drinking water, raw rice and cooking of rice as arsenic
305 exposure routes in three contrasting areas of West Bengal, India. *Environmental geochemistry*
306 *and health* **32**, 463-477 (2010).
- 307 4. WHO, Guidelines for drinking-water quality. *WHO chronicle* **38**, 104-108 (2011).
- 308 5. P. Smedley, D. Kinniburgh, A review of the source, behaviour and distribution of arsenic in
309 natural waters. *Applied geochemistry* **17**, 517-568 (2002).
- 310 6. R. Nickson *et al.*, Arsenic poisoning of Bangladesh groundwater. *Nature* **395**, 338-338 (1998).
- 311 7. J. McArthur, P. Ravenscroft, S. Safiulla, M. Thirlwall, Arsenic in groundwater: testing pollution
312 mechanisms for sedimentary aquifers in Bangladesh. *Water Resources Research* **37**, 109-117
313 (2001).
- 314 8. M. Berg *et al.*, Hydrological and sedimentary controls leading to arsenic contamination of
315 groundwater in the Hanoi area, Vietnam: the impact of iron-arsenic ratios, peat, river bank
316 deposits, and excessive groundwater abstraction. *Chemical Geology* **249**, 91-112 (2008).
- 317 9. S. Fendorf, H. A. Michael, A. van Geen, Spatial and temporal variations of groundwater arsenic
318 in South and Southeast Asia. *Science* **328**, 1123-1127 (2010).
- 319 10. M. I. Litter *et al.*, Arsenic in Argentina: Occurrence, human health, legislation and
320 determination. *Science of The Total Environment*, (2019).
- 321 11. Y. Zheng, S. V. Flanagan, The case for universal screening of private well water quality in the US
322 and testing requirements to achieve it: evidence from arsenic. *Environmental health*
323 *perspectives* **125**, 085002 (2017).
- 324 12. P. Bhattacharya, D. Chatterjee, G. Jacks, Occurrence of Arsenic-contaminated Groundwater in
325 Alluvial Aquifers from Delta Plains, Eastern India: Options for Safe Drinking Water Supply.
326 *International Journal of Water Resources Development* **13**, 79-92 (1997).

- 327 13. A. Van Geen *et al.*, Spatial variability of arsenic in 6000 tube wells in a 25 km² area of
328 Bangladesh. *Water Resources Research* **39**, (2003).
- 329 14. H. B. Nicolli, J. M. Suriano, M. A. G. Peral, L. H. Ferpozzi, O. A. Baleani, Groundwater
330 contamination with arsenic and other trace elements in an area of the Pampa, Province of
331 Córdoba, Argentina. *Environmental Geology and Water Sciences* **14**, 3-16 (1989).
- 332 15. P. Smedley, H. Nicolli, D. Macdonald, A. Barros, J. Tullio, Hydrogeochemistry of arsenic and
333 other inorganic constituents in groundwaters from La Pampa, Argentina. *Applied Geochemistry*
334 **17**, 259-284 (2002).
- 335 16. M. Blarasin, A. Cabrera, E. Matteoda, paper presented at the XXXIII IAH -. 7^o ALHSUD Congress,
336 Zacatecas, Mexico, 11 - 15 October 2004 2004.
- 337 17. M. Auge, G. E. Viale, L. Sierra, in *VIII Congreso Argentino de Hidrogeología: Aguas subterráneas*
338 *recurso estratégico*. (2013), vol. 2, pp. 58-63.
- 339 18. M. Berg *et al.*, Arsenic contamination of groundwater and drinking water in Vietnam: a human
340 health threat. *Environmental Science & Technology* **35**, 2621-2626 (2001).
- 341 19. J. Buschmann, M. Berg, C. Stengel, M. L. Sampson, Arsenic and manganese contamination of
342 drinking water resources in Cambodia: coincidence of risk areas with low relief topography.
343 *Environmental science & technology* **41**, 2146-2152 (2007).
- 344 20. P. Smedley, M. Zhang, G. Zhang, Z. Luo, Mobilisation of arsenic and other trace elements in
345 fluviolacustrine aquifers of the Huhhot Basin, Inner Mongolia. *Applied Geochemistry* **18**, 1453-
346 1477 (2003).
- 347 21. L. Rodríguez-Lado *et al.*, Groundwater arsenic contamination throughout China. *Science* **341**,
348 866-868 (2013).
- 349 22. Y. Zhou *et al.*, Distribution of groundwater arsenic in Xinjiang, PR China. *Applied geochemistry*
350 **77**, 116-125 (2017).
- 351 23. D. Chatterjee, R. Roy, B. Basu, Riddle of arsenic in groundwater of Bengal Delta Plain—role of
352 non-inland source and redox traps. *Environmental Geology* **49**, 188-206 (2005).
- 353 24. B. Nath, D. Stüben, S. B. Mallik, D. Chatterjee, L. Charlet, Mobility of arsenic in West Bengal
354 aquifers conducting low and high groundwater arsenic. Part I: Comparative hydrochemical and
355 hydrogeological characteristics. *Applied Geochemistry* **23**, 977-995 (2008).
- 356 25. B. A. Shah, Arsenic-contaminated groundwater in Holocene sediments from parts of middle
357 Ganga plain, Uttar Pradesh, India. *Current Science(Bangalore)* **98**, 1359-1365 (2010).
- 358 26. B. Planer-Friedrich, *Hydrogeological and hydrochemical investigations in the Rioverde basin,*
359 *Mexico*. (Verlag nicht ermittelbar, 2000).
- 360 27. M. T. Alarcón-Herrera *et al.*, Co-occurrence of arsenic and fluoride in groundwater of semi-arid
361 regions in Latin America: Genesis, mobility and remediation. *Journal of Hazardous Materials*
362 **262**, 960-969 (2013).
- 363 28. R. Nickson, J. McArthur, B. Shrestha, T. Kyaw-Myint, D. Lowry, Arsenic and other drinking
364 water quality issues, Muzaffargarh District, Pakistan. *Applied Geochemistry* **20**, 55-68 (2005).
- 365 29. J. E. Podgorski *et al.*, Extensive arsenic contamination in high-pH unconfined aquifers in the
366 Indus Valley. *Science Advances* **3**, (2017).
- 367 30. J. D. Ayotte, M. G. Nielsen, G. R. Robinson Jr, R. B. Moore, Relation of arsenic, iron, and
368 manganese in ground water to aquifer type, bedrock lithochemistry, and land use in the
369 New England Coastal Basins. *Water Resources Investigations Report* **99**, 4162 (1999).
- 370 31. J. D. Ayotte, L. Medalie, S. L. Qi, L. C. Backer, B. T. Nolan, Estimating the high-arsenic domestic-
371 well population in the conterminous United States. *Environ. Sci. Technol.* **51**, 12443-12454
372 (2017).
- 373 32. M. Berg *et al.*, Magnitude of arsenic pollution in the Mekong and Red River Deltas—Cambodia
374 and Vietnam. *Science of the Total Environment* **372**, 413-425 (2007).
- 375 33. L. H. Winkel *et al.*, Arsenic pollution of groundwater in Vietnam exacerbated by deep aquifer
376 exploitation for more than a century. *Proceedings of the National Academy of Sciences* **108**,
377 1246-1251 (2011).
- 378 34. A. Bretzler *et al.*, Groundwater arsenic contamination in Burkina Faso, West Africa: Predicting
379 and verifying regions at risk. *Science of the Total Environment* **584**, 958-970 (2017).

- 380 35. Q. Zhang *et al.*, Coupling predicted model of arsenic in groundwater with endemic arsenism
381 occurrence in Shanxi Province, Northern China. *Journal of hazardous materials* **262**, 1147-1153
382 (2013).
- 383 36. S. Bindal, C. K. Singh, Predicting groundwater arsenic contamination: regions at risk in highest
384 populated state of India. *Water Res* **159**, 65-76 (2019).
- 385 37. L. Winkel, M. Berg, M. Amini, S. J. Hug, C. A. Johnson, Predicting groundwater arsenic
386 contamination in Southeast Asia from surface parameters. *Nat. Geosci.* **1**, 536-542 (2008).
- 387 38. Q. Yang, H. B. Jung, R. G. Marvinney, C. W. Culbertson, Y. Zheng, Can arsenic occurrence rates
388 in bedrock aquifers be predicted? *Environmental science & technology* **46**, 2080-2087 (2012).
- 389 39. N. Yang, L. H. Winkel, K. H. Johannesson, Predicting geogenic arsenic contamination in shallow
390 groundwater of South Louisiana, United States. *Environmental science & technology* **48**, 5660-
391 5666 (2014).
- 392 40. M. Amini *et al.*, Statistical modeling of global geogenic arsenic contamination in groundwater.
393 *Environmental science & technology* **42**, 3669-3675 (2008).
- 394 41. P. Ravenscroft, H. Brammer, K. Richards, *Arsenic pollution: a global synthesis*. (John Wiley &
395 Sons, 2009), vol. 28.
- 396 42. S. Murcott, *Arsenic contamination in the world*. (IWA publishing, 2012).
- 397 43. P. Ravenscroft, Predicting the global extent of arsenic pollution of groundwater and its
398 potential impact on human health. *Unpublished report prepared for UNICEF, December,*
399 (2007).
- 400 44. Y. Zheng, Lessons learned from arsenic mitigation among private well households. *Current*
401 *environmental health reports* **4**, 373-382 (2017).
- 402 45. J. Podgorski, M. Berg, Podgorski_and_Berg_2020. *ERIC/open* (2020); doi.org/10.25678/0001ZT.
- 403 46. L. Breiman, Random forests. *Machine learning* **45**, 5-32 (2001).
- 404 47. R Core Team, R Foundation for Statistical Computing, Ed. (Vienna, Austria, 2014).
- 405 48. T. T. Hastie, Robert; Friedman, Jerome, *The Elements of Statistical Learning (2nd ed.)*.
406 (Springer, 2008).
- 407 49. T. K. Ho, in *Document analysis and recognition, 1995., proceedings of the third international*
408 *conference on*. (IEEE, 1995), vol. 1, pp. 278-282.
- 409 50. M. L. McHugh, Interrater reliability: the kappa statistic. *Biochemia medica: Biochemia medica*
410 **22**, 276-282 (2012).
- 411 51. T. Fawcett, An introduction to ROC analysis. *Pattern recognition letters* **27**, 861-874 (2006).
- 412 52. J. Gao, NASA Socioeconomic Data and Applications Center (SEDAC), Ed. (Palisades, NY, 2019).
- 413 53. JMP. (WHO/UNICEF Joint Monitoring Program (JMP), 2019), vol. 2019.
- 414 54. M. A. Friedl *et al.*, MODIS Collection 5 global land cover: Algorithm refinements and
415 characterization of new datasets. *Remote sensing of Environment* **114**, 168-182 (2010).
- 416 55. FAO, Food and Agriculture Organization of the United Nations (FAO), Ed. (2016).
- 417 56. R. E. Broshears, M. A. Akbari, M. P. Chornack, D. K. Mueller, B. C. Ruddy, "Inventory of ground-
418 water resources in the Kabul Basin, Afghanistan," (U. S. Geological Survey, 2005).
- 419 57. UNHCR. (UNHCR, 2019), vol. 2019.
- 420 58. M. E. Zabala, M. Manzano, L. Vives, Assessment of processes controlling the regional
421 distribution of fluoride and arsenic in groundwater of the Pampeano Aquifer in the Del Azul
422 Creek basin (Argentina). *Journal of hydrology* **541**, 1067-1087 (2016).
- 423 59. M. E. Morgada, M. Mateu, J. Bundschuh, M. I. Litter, Arsenic in the Iberoamerican region. The
424 IBEROARSEN Network and a possible economic solution for arsenic removal in isolated rural
425 zones. *e-Terra* **5**, 1-11 (2008).
- 426 60. K. Ivkovic, K. Watkins, R. Cresswell, J. Bauld, A groundwater quality assessment of the
427 fractured rock aquifers of the Piccadilly Valley, South Australia. (1998).
- 428 61. J. Fitzgerald *et al.*, Groundwater quality and environmental health implications. *Anangu*
429 *Pitjantjara Lands, South Australia, A Report from Bureau of Rural Sciences*, 1-30 (1999).
- 430 62. S. Clohessy, "Perth Shallow Groundwater Systems Investigation: Lake Gwelup,"
431 *Hydrogeological record series* (Department of Water, Perth, 2012).

- 432 63. R. M. Larsen, *A Groundwater Quality Assessment of the Jandakot Mound, Swan Coastal Plain, Western Australia*. (Australian Geological Survey Organisation, 1998).
- 433
- 434 64. E. E. A. EEA. (2019).
- 435 65. D. Kinniburgh, P. Smedley, Arsenic contamination of groundwater in Bangladesh. (2001).
- 436 66. *DWA Maun Groundwater Development Project: Phase 2, Resources assessment and wellfield development: Final report*, (2004).
- 437
- 438 67. CPRM, Geological Survey of Brazil, Ed. (2017), vol. 2017.
- 439 68. Ministry of Rural Development of Cambodia. (2015).
- 440 69. E. K. Read *et al.*, Water quality data for national-scale aquatic research: The Water Quality Portal. *Water Resources Research* **53**, 1735-1745 (2017).
- 441
- 442 70. D. R. Boyle, W. A. Spirito, S. W. Adcock, "Groundwater hydrogeochemical survey of central New Brunswick," (1996).
- 443
- 444 71. C. Reimann, K. Bjorvatn, R. Tekle-Haimanot, Z. Melako, U. Siewers, Drinking water quality, Rift Valley, Ethiopia. *Norges geologiske undersøkelse, Report* **2002**, 132 (2002).
- 445
- 446 72. A. Bretzler *et al.*, Groundwater origin and flow dynamics in active rift systems—A multi-isotope approach in the Main Ethiopian Rift. *Journal of hydrology* **402**, 274-289 (2011).
- 447
- 448 73. T. Rango, G. Bianchini, L. Beccaluva, R. Tassinari, Geochemistry and water quality assessment of central Main Ethiopian Rift natural waters with emphasis on source and occurrence of fluoride and arsenic. *Journal of African Earth Sciences* **57**, 479-491 (2010).
- 449
- 450
- 451 74. B. Kortatsi *et al.*, Reconnaissance survey of arsenic concentration in ground-water in south-eastern Ghana. *West African Journal of Applied Ecology* **13**, 16-26 (2008).
- 452
- 453 75. P. L. Smedley, Arsenic in rural groundwater in Ghana: part special issue: hydrogeochemical studies in sub-Saharan Africa. *Journal of African Earth Sciences* **22**, 459-470 (1996).
- 454
- 455 76. I. A. Katsoyiannis, S. J. Hug, A. Ammann, A. Zikoudi, C. Hatziliontos, Arsenic speciation and uranium concentrations in drinking water supply wells in Northern Greece: correlations with redox indicative parameters and implications for groundwater treatment. *Science of the Total Environment* **383**, 128-140 (2007).
- 456
- 457
- 458
- 459 77. H. A. Rowland *et al.*, Geochemistry and arsenic behaviour in groundwater resources of the Pannonian Basin (Hungary and Romania). *Applied Geochemistry* **26**, 1-17 (2011).
- 460
- 461 78. S. Chandra, S. Ahmed, E. Nagaiah, S. K. Singh, P. Chandra, Geophysical exploration for lithological control of arsenic contamination in groundwater in Middle Ganga Plains, India. *Physics and Chemistry of the Earth, Parts A/B/C* **36**, 1353-1362 (2011).
- 462
- 463
- 464 79. T. Ghosh, R. Kanchan, Geoenvironmental appraisal of groundwater quality in Bengal alluvial tract, India: a geochemical and statistical approach. *Environmental earth sciences* **72**, 2475-2488 (2014).
- 465
- 466
- 467 80. A. Mukherjee *et al.*, Controls on high and low groundwater arsenic on the opposite banks of the lower reaches of River Ganges, Bengal basin, India. *Science of the Total Environment* **645**, 1371-1387 (2018).
- 468
- 469
- 470 81. J. McArthur *et al.*, How paleosols influence groundwater flow and arsenic pollution: a model from the Bengal Basin and its worldwide implication. *Water Resources Research* **44**, (2008).
- 471
- 472 82. A. Mukherjee, A. E. Fryar, H. D. Rowe, Regional-scale stable isotopic signatures of recharge and deep groundwater in the arsenic affected areas of West Bengal, India. *Journal of Hydrology* **334**, 151-161 (2007).
- 473
- 474
- 475 83. V. S. Chauhan, R. Nickson, D. Chauhan, L. Iyengar, N. Sankararamakrishnan, Ground water geochemistry of Ballia district, Uttar Pradesh, India and mechanism of arsenic release. *Chemosphere* **75**, 83-91 (2009).
- 476
- 477
- 478 84. D. Saha, S. Sahu, A decade of investigations on groundwater arsenic contamination in Middle Ganga Plain, India. *Environmental geochemistry and health* **38**, 315-337 (2016).
- 479
- 480 85. D. P. Shukla, C. Dubey, N. P. Singh, M. Tajbakhsh, M. Chaudhry, Sources and controls of Arsenic contamination in groundwater of Rajnandgaon and Kanker District, Chattisgarh Central India. *Journal of Hydrology* **395**, 49-66 (2010).
- 481
- 482

- 483 86. J. P. Maity *et al.*, Arsenic-enriched groundwaters of India, Bangladesh and Taiwan—
484 Comparison of hydrochemical characteristics and mobility constraints. *Journal of*
485 *Environmental Science and Health, Part A* **46**, 1163-1176 (2011).
- 486 87. S. Kar *et al.*, Arsenic-enriched aquifers: occurrences and mobilization of arsenic in groundwater
487 of Ganges Delta Plain, Barasat, West Bengal, India. *Applied Geochemistry* **25**, 1805-1814
488 (2010).
- 489 88. S. Hazarika, B. Bhuyan, Fluoride, arsenic and iron content of groundwater around six selected
490 tea gardens of Lakhimpur District, Assam, India. *Arch Appl Sci Res* **5**, 57-61 (2013).
- 491 89. B. Nath *et al.*, Hydrochemistry of arsenic-enriched aquifer from rural West Bengal, India: a
492 study of the arsenic exposure and mitigation option. *Water, air, and soil pollution* **190**, 95-113
493 (2008).
- 494 90. R. A. Olea, N. J. Raju, J. J. Egozcue, V. Pawlowsky-Glahn, S. Singh, Advancements in
495 hydrochemistry mapping: methods and application to groundwater arsenic and iron
496 concentrations in Varanasi, Uttar Pradesh, India. *Stochastic environmental research and risk*
497 *assessment* **32**, 241-259 (2018).
- 498 91. M. Kumar, A. Ramanathan, M. M. Rahman, R. Naidu, Concentrations of inorganic arsenic in
499 groundwater, agricultural soils and subsurface sediments from the middle Gangetic plain of
500 Bihar, India. *Science of the Total Environment* **573**, 1103-1114 (2016).
- 501 92. S. Chidambaram *et al.*, A study on the arsenic concentration in groundwater of a coastal
502 aquifer in south-east India: an integrated approach. *Environment, Development and*
503 *Sustainability* **19**, 1015-1040 (2017).
- 504 93. S. Ghosh, P. Sar, Identification and characterization of metabolic properties of bacterial
505 populations recovered from arsenic contaminated ground water of North East India (Assam).
506 *Water research* **47**, 6992-7005 (2013).
- 507 94. S. Sharma, J. Kaur, A. K. Nagpal, I. Kaur, Quantitative assessment of possible human health risk
508 associated with consumption of arsenic contaminated groundwater and wheat grains from
509 Ropar Wetland and its environs. *Environmental monitoring and assessment* **188**, 506 (2016).
- 510 95. B. A. Shah, Role of Quaternary stratigraphy on arsenic-contaminated groundwater from parts
511 of Barak Valley, Assam, North–East India. *Environmental earth sciences* **66**, 2491-2501 (2012).
- 512 96. B. A. Shah, Role of Quaternary stratigraphy on arsenic-contaminated groundwater from parts
513 of Middle Ganga Plain, UP–Bihar, India. *Environmental geology* **53**, 1553-1561 (2008).
- 514 97. B. A. Shah, Status of groundwater arsenic pollution of Mirzapur district in Holocene aquifers
515 from parts of the Middle Ganga Plain, India. *Environmental earth sciences* **73**, 1505-1514
516 (2015).
- 517 98. L. Sailo, C. Mahanta, Arsenic mobilization in the Brahmaputra plains of Assam: groundwater
518 and sedimentary controls. *Environmental monitoring and assessment* **186**, 6805-6820 (2014).
- 519 99. D. Paul, S. K. Kazy, A. K. Gupta, T. Pal, P. Sar, Diversity, metabolic properties and arsenic
520 mobilization potential of indigenous bacteria in arsenic contaminated groundwater of West
521 Bengal, India. *PloS one* **10**, e0118735 (2015).
- 522 100. Catholic Relief Services, Meulaboh Water Quality Laboratory in Aceh, Ed. (2007).
- 523 101. UNEP, "Water Quality, 2005 State of the UNEP GEMS/Water Global Network and Annual
524 Report," (2005).
- 525 102. M. Pritchard, T. Mkandawire, J. O'neill, Assessment of groundwater quality in shallow wells
526 within the southern districts of Malawi. *Physics and Chemistry of the Earth, Parts A/B/C* **33**,
527 812-823 (2008).
- 528 103. INCA, *Arsénico y fluoruro en agua: riesgos y perspectivas desde la sociedad civil y la academia*
529 *en México*. (Mexico, 2018).
- 530 104. A. Van Geen *et al.*, Confirmation of elevated arsenic levels in groundwater of Myanmar.
531 *Science of the Total Environment* **478**, 21-24 (2014).
- 532 105. B. R. Shrestha, J. W. Whitney, K. B. Shrestha, "The State of Arsenic in Nepal-2003," (
533 Kathmandu, Nepal, 2004).

- 534 106. B. Frengstad, A. K. M. Skrede, D. Banks, J. R. Krog, U. Siewers, The chemistry of Norwegian
535 groundwaters: III. The distribution of trace elements in 476 crystalline bedrock groundwaters,
536 as analysed by ICP-MS techniques. *Science of the Total environment* **246**, 21-40 (2000).
- 537 107. P. de Caritat, S. Danilova, C. Reimann, G. Storrø, Groundwater composition near the nickel—
538 copper smelting industry on the Kola Peninsula, central Barents Region (NW Russia and NE
539 Norway). *Journal of Hydrology* **208**, 92-107 (1998).
- 540 108. C. M. de Meyer *et al.*, Arsenic, manganese and aluminum contamination in groundwater
541 resources of Western Amazonia (Peru). *Science of the Total Environment* **607**, 1437-1450
542 (2017).
- 543 109. L. McCaffrey, J. Willis, *Distribution of fluoride-rich groundwater in the eastern and Mogwase*
544 *regions of the Northern and North-West Provinces*. (Water Research Commission Pretoria,
545 2001).
- 546 110. Geological Survey of Sweden. (2007).
- 547 111. M. Haldimann, E. Pfammatter, P.-M. Venetz, P. Studer, V. Dudler, Occurrence of arsenic in
548 drinking water of the canton of Valais. Part I: Overview of arsenic concentration and
549 geographic distribution. *Mitteilungen aus Lebensmitteluntersuchung und Hygiene* **96**, 89-105
550 (2005).
- 551 112. P. Smedley *et al.*, Fluoride in groundwater from high-fluoride areas of Ghana and Tanzania.
552 (2002).
- 553 113. J. Buschmann *et al.*, Contamination of drinking water resources in the Mekong delta
554 floodplains: Arsenic and other trace metals pose serious health risks to population.
555 *Environment International* **34**, 756-764 (2008).
- 556 114. A. Trabucco, R. Zomer, Global soil water balance geospatial database. *CGIAR Consortium for*
557 *Spatial Information, Published online, available from the CGIAR-CSI GeoPortal at:*
558 *http://www.cgiar-csi.org (last access: January 2013),* (2010).
- 559 115. A. Trabucco, R. J. Zomer, Global aridity index (global-aridity) and global potential evapo-
560 transpiration (global-PET) geospatial database. *CGIAR Consortium for Spatial Information,*
561 (2009).
- 562 116. R. J. Hijmans, S. E. Cameron, J. L. Parra, P. G. Jones, A. Jarvis, Very high resolution interpolated
563 climate surfaces for global land areas. *International journal of climatology* **25**, 1965-1978
564 (2005).
- 565 117. S. E. Fick, R. J. Hijmans, WorldClim 2: new 1-km spatial resolution climate surfaces for global
566 land areas. *International Journal of Climatology*, (2017).
- 567 118. J. Hartmann, N. Moosdorf, The new global lithological map database GLIM: A representation of
568 rock properties at the Earth surface. *Geochemistry, Geophysics, Geosystems* **13**, (2012).
- 569 119. C. E. R. S. C. U.S. Geological Survey - Energy Resources Program. (Department of the
570 Interior/USGS).
- 571 120. T. Hengl *et al.*, SoilGrids250m: Global gridded soil information based on machine learning. *PLoS*
572 *one* **12**, e0169748 (2017).
- 573 121. C. W. Ross *et al.*, HYSOGs250m, global gridded hydrologic soil groups for curve-number-based
574 runoff modeling. *Scientific data* **5**, 180091 (2018).
- 575 122. J. Pelletier *et al.*, Global 1-km gridded thickness of soil, regolith, and sedimentary deposit
576 layers. *ORNL DAAC*, (2016).
- 577 123. USGS. (U.S. Geological Survey, EROS Data Cent. Sioux Falls, SD, 1996).
- 578 124. T. Hengl. (Zenodo, 2018).
- 579 125. Y. Fan, H. Li, G. Miguez-Macho, Global patterns of groundwater table depth. *Science* **339**, 940-
580 943 (2013).

581